

## Supplementary Information

Worobey *et al.* “Direct Evidence of Extensive Diversity of HIV-1 in Kinshasa by 1960”

### This file contains

Supplementary text

Supplementary Table S1

Supplementary Figs. S1, S2, and S3 and figure legends

**Authenticity of the data.** Several lines of evidence argue strongly that the sequences recovered from the 1960 sample are genuine and not the result of contamination:

- (i) The work was performed in laboratories dedicated to working with ancient/degraded samples, under strict controls for contamination, including physical separation of pre- and post-PCR work and numerous extraction and PCR blanks that were consistently negative;
- (ii) The results were reproducible across multiple extractions followed by RT-PCR, cloning, and sequencing;
- (iii) The results were independently replicated in SMW’s laboratory (Northwestern University) using three blinded sub-samples of lymph node tissue; the correct sample was identified as the positive by RT-PCR and verified by cloning and sequencing (**Table S1**);
- (iv) 12 short RNA fragments, covering the *gag*, *pol*, and *env* genes, were successfully amplified by RT-PCR. No HIV-1 amplicons larger than 126 nucleotides (including the longer sets used by Zhu *et al.* to amplify ZR59) were recoverable from the sample, as predicted for a Bouin’s-fixed paraffin-embedded specimen;

- (v) Similarly, the qPCR assay using a human *B2M* RNA marker confirmed that amplifiable RNA was present in the sample but that its quality was relatively low, as expected for paraffin-embedded tissues;
- (vi) As observed with ZR59, the relatively short branch length of the HIV-1 sequences from 1960 is consistent with it being an early strain;
- (vii) The sample originated from a patient (an adult with apparent lymph node abnormalities) and a tissue type (lymph node) where HIV-1 RNA could plausibly be present and detectable; and
- (viii) All of the paraffin-embedded tissues yielded sequences with reasonable topological and branch length properties (*i.e.*, subtype B sequences from the Canadian sample, and subtype D and A sequences from the three linked to the Democratic Republic of the Congo, each closely related to previously published sequences of Congolese origin).

**Table S1.** Primers and sequence-confirmed amplification results.

Primer	Primer Sequence 5'-3'	Frag- ment #	DRC 60	BE81	BE85	CAN97
HIVG1 (F) <sup>a</sup>	ACCCACCTATCCCAGTAGGAGAAAT	1	(+) <sup>c</sup>	(+)	(+)	(+)
HIVG2 (R)	GGTCCTTGTCTTATGTCCAGAATGC					
Pol3290F <sup>b</sup>	GCCAGAAAAAGACAGCTGGACTGTCAA	2	(+) <sup>c</sup>	(+)	(+)	(+)
Pol3415R	AGAGCTTTGGYTCCCCTAAGG					
env6371F	CACCACTCTATTTTGTGCATCAG	3	(+)	(+)	(+)	(-)
env6442R	GCATGTGTAGCCCAGACATTAT					
env6445F	GTGTACCCACAGACCCCAAC	4	(+)	(+)	(+)	(-)
env6542R	CTCATGCATTTGTTCTACCATGT					
env6828F	ACACAGGCTTGTCCAAAGGT	5	(+) <sup>c</sup>	(-)	(-)	(-)
env6890R	ACCAGCTGGGGCACAATAAT					
env7468F	CACTCCCATGCAGAATAAAACA	6	(+) <sup>c</sup>	(+)	(+)	(-)
env7535R	AGGGGCATACATTGCTTGTC					
env7717F <sup>d</sup>	CCACCAAGGCAAAGAGAAGA	7	(-)	(+)	(-)	(+)
env7796R	TCCCAAGAACCCAAGGAAC					
env7899F	ATAGAGGCGCAACAGCATCT	8	(-)	(-)	(-)	(-)
env7977R	TTTCCACAGCCAGGACTCTT					
env8047F	TGCCCTGGAACCTAGTTGG	9	(+) <sup>c</sup>	(+)	(+)	(-)
env8112R	CCATCCAGGTCATGTTCTCC					
env8423F	CGAAGAAGAAGGTGGAGAGC	10	(+) <sup>c</sup>	(+)	(-)	(-)
env8498R	GTCCCAGGCAAGTGCTAAGA					
env7717F <sup>d,e</sup>	CCACCAAGGCAAAGAGAAGA	11	(+)	N/D <sup>f</sup>	N/D	N/D
env7805R	TCCTGCTGCTCCTAAGAACC					
env7771F <sup>d,e</sup>	GAGCTGTCTTCCTTGGGTTCT	12	(+)	N/D	N/D	N/D
env7846R	GCCTGTACCGTCAGCGTTA					
env7835F <sup>d,e</sup>	GACGGTACAGGCCAGACAAT	13	(+)	N/D	N/D	N/D
env7937R	CCAGACCGTGAGTTTCAACA					
env7890F <sup>d,e</sup>	CTGAGGGCTATAGAGGCTCAAC	14	(+) <sup>c</sup>	N/D	N/D	N/D
env7960R	CTTGCCTGGAGCTGTTTAATG					

<sup>a</sup>Previously published [Boni, J. & Schupbach, J. *J. Virol. Methods* **42**, 309 (1993)].

<sup>b</sup>The remaining primer names correspond to the primer's 5' base location in the HXB2 reference strain.

<sup>c</sup>Independently replicated in SMW's laboratory, Northwestern University.

<sup>d</sup>Overlaps ZR59 sequence.

<sup>e</sup>Designed specifically for HIV-1 group M subtype A (after DRC60 was determined to be A-like).

<sup>f</sup>N/D, no data: amplification not attempted.

**Fig. S1.** The 50% majority-rule consensus tree summarizing the Bayesian MCMC phylogenetic analysis conducted using MrBayes v3.1.2. Branch lengths (substitutions/site) represent the average for that branch in tree sample. The posterior probability for each node is indicated. DRC60 (red), ZR59 (black), and the three control sequences recovered from paraffin-embedded specimens (gray) are depicted with bold branches. DRC60 and the two positive controls from the DRC form monophyletic clades with other DRC sequences, while the Canadian control falls as expected within subtype B. Each control sequence has a reasonably long terminal branch length consistent with its date of sampling. The subtype A/A1 ancestral node used for the regression line dating analysis is indicated. The country of sampling and additional details about the reference strains used can be found at <http://hiv.lanl.gov>.

**Fig. S2.** The maximum clade credibility tree from the relaxed clock Bayesian MCMC analysis conducted using BEAST v1.4.7. This is the same tree as shown in **Fig. 2**, but in rectangular form and with sequence names indicated. The country of sampling and additional details about the reference strains used can be found at <http://hiv.lanl.gov>.

**Fig. S3.** Bayesian skyline plot of HIV-1 group M. The plot begins at the median posterior TMRCA (1908). The bold line traces the inferred median effective population size over time with the 95% HPD shaded in blue.

Fig. S1

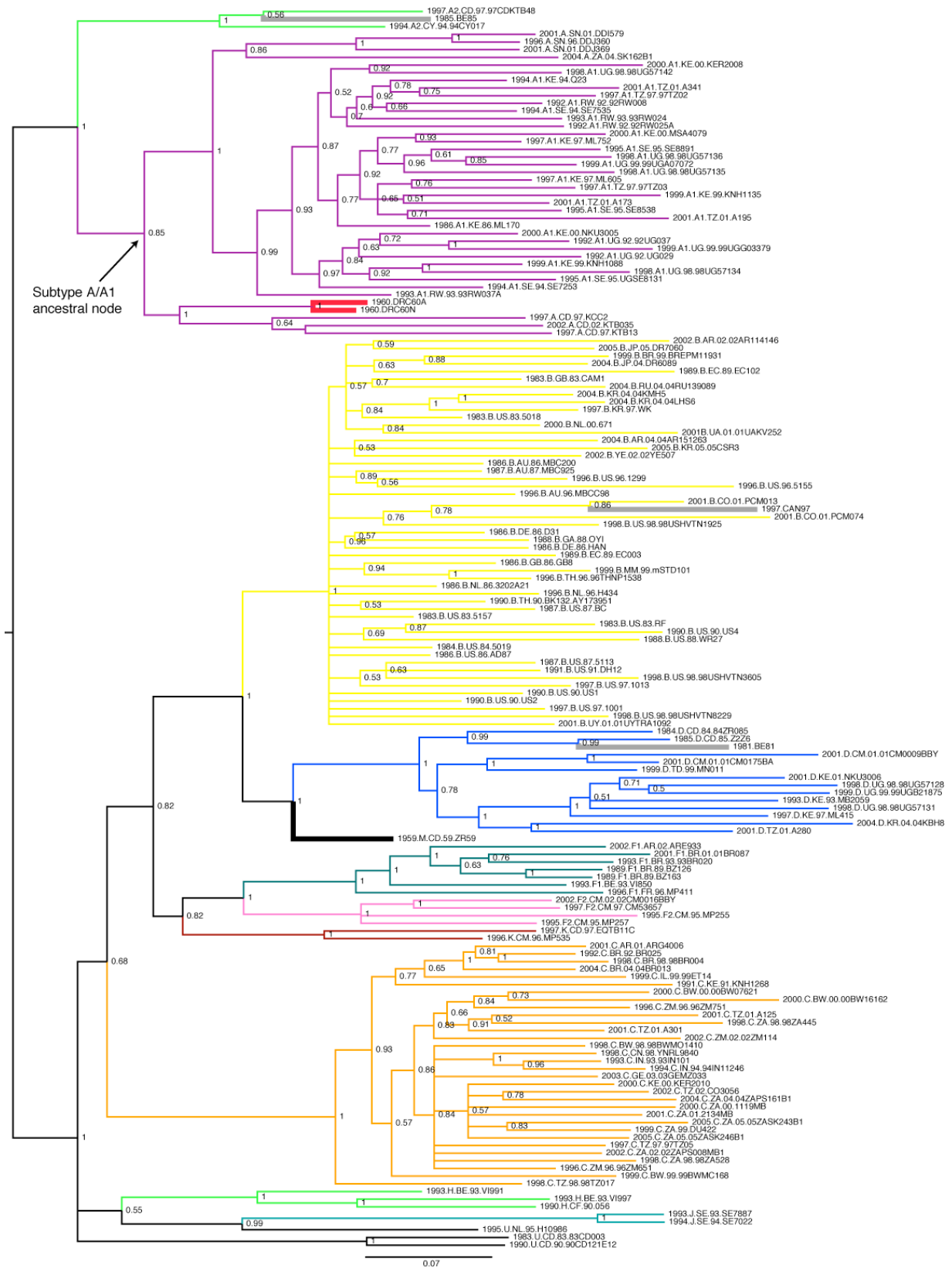


Fig. S2



**Fig. S3**

